# Linear models for classification

$$\underset{\substack{\text{input,}\\ D \text{ dims}}}{\vec{x}} \;\to\; \underset{\substack{\uparrow\\ \text{discrete classes}}}{C_k} \;,\; k = 1, \ldots, K$$

Input space divided into <u>decision regions</u> by <u>decision surfaces</u>.

For ex., $\underset{\substack{\text{target}\\ \text{variables}}}{\vec{t}} = \underbrace{(0, 1, 0, 0, 0)}$

$K = 5$ classes, target variable indicates class 2

Classification approaches:

(1) Discriminant function: directly assigns $\vec{x}$ to a class, e.g. for 2 classes

$$y(\vec{x}) \geq 0 \iff C_1, \quad y(\vec{x}) < 0 \iff C_2.$$

(2) Probabilistic approach:

use $p(C_k | \vec{x}) = \underbrace{\dfrac{p(\vec{x} | C_k)\, p(C_k)}{p(\vec{x})}}_{\text{Bayesian framework}}$

Previously, we focused on

$$y(\vec{x}) = \vec{w}^T \vec{\varphi}(\vec{x}) + w_0 \;\overset{\text{explicit bias prm}}{.}$$

Now we will consider $y(\vec{x}) = f(\vec{w}^T \vec{x} + w_0)$

or, more generally, $\leftarrow$ ↑

$$y(\vec{x}) = f(\vec{w}^T \vec{\varphi}(\vec{x}) + w_0) \qquad \substack{\text{non-linear}\\ \text{activation } \underline{f'n}}$$

Decision surfaces are given by

$$\vec{w}^T \vec{x} + w_0 = \text{const} \quad , \quad \text{linear f's of } \vec{x}$$

Discriminant functions

① Two classes

Consider $y(\vec{x}) = \vec{w}^T \vec{x} + w_0$ [linear discriminant]

$$\begin{cases} y(\vec{x}) \geqslant 0 \Rightarrow C_1 \\ y(\vec{x}) < 0 \Rightarrow C_2 \end{cases}$$

$y(\vec{x}) = 0 \Leftarrow$ decision boundary (DB)

Consider $\vec{x}_A, \vec{x}_B \in DB$

$$y(\vec{x}_A) = y(\vec{x}_B) = 0 \Rightarrow \vec{w}^T \underbrace{(\vec{x}_A - \vec{x}_B)}_{\substack{\text{lies on } (D-1) \\ \text{dim'l DB}}} = 0$$

So, $\vec{w} \perp DB \Rightarrow \vec{n} = \dfrac{\vec{w}}{\|\vec{w}\|}$ is a unit vector $\perp DB$
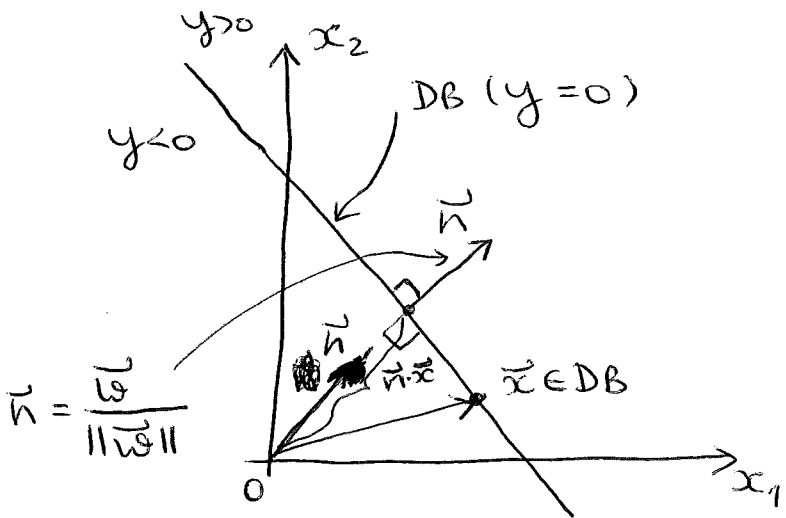
Similarly, if $\vec{x} \in DB$,

$$y(\vec{x}) = \vec{w}^T \vec{x} + w_0 = 0 \quad , \quad \text{or}$$

$$\underbrace{\frac{\vec{w}^T \vec{x}}{\|\vec{w}\|}}_{} = -\frac{w_0}{\|\vec{w}\|}$$

" $\vec{n} \cdot \vec{x}$ = normal distance from the origin to DB

So, $w_0$ determines the location of DB

Moreover, for any $\vec{x}$

$$\vec{x} = \vec{x}_{\parallel DB} + \underbrace{r\vec{n}}_{\perp DB}$$

Then

$$\underbrace{\vec{\omega}^T \vec{x} + \omega_0}_{y(\vec{x})} = \vec{\omega}^T \vec{x}_{\parallel DB} + r \frac{\overbrace{\vec{\omega}^T \cdot \vec{\omega}}^{\|\vec{\omega}\|^2}}{\|\vec{\omega}\|} + \omega_0, \text{ or}$$

$$y(\vec{x}) = r\|\vec{\omega}\| \implies r = \frac{y(\vec{x})}{\|\vec{\omega}\|}.$$

↑ perpendicular distance
• $\vec{x}$ to DB
from

Here, we used

$$\vec{\omega}^T \cdot \vec{x}_{\parallel DB} + \omega_0 = 0$$

Indeed,



- 3 -

Then

$$-\vec{n}\,\frac{w_0}{\|\vec{w}\|} + \vec{x}_2 = \vec{x}$$

$$\underbrace{\vec{x}_{2,\|DB}}_{} + \underbrace{\vec{x}_{2,\perp DB}}_{\text{"}r\vec{n}}$$

$$\vec{w}^T\cdot\vec{x} = -\underbrace{\frac{\vec{w}^T\cdot\vec{w}}{\|\vec{w}\|^2}}_{\text{"}1}w_0 + \underbrace{\vec{w}^T\cdot\vec{x}_{2,\|DB}}_{\text{"}0} + r\underbrace{\frac{\vec{w}^T\cdot\vec{w}}{\|\vec{w}\|}}_{\|\vec{w}\|}\,,\ \text{or}$$

$$\underbrace{\vec{w}^T\cdot\vec{x} + w_0}_{y(\vec{x})} = r\,\|\vec{w}\|\quad,\ \text{as before.}$$

② <u>Multiple classes</u>

Consider $K > 2$ classes.
Difficulties in generalizing from the $K=2$ case:



majority vote among discriminant $f$'s

Rather, consider a single $K$-class discriminant :

$$y_K(\vec{x}) = \vec{w}_K^T\cdot\vec{x} + w_{K,0}$$

Assign a point to $C_K$ if $y_K(\vec{x}) > y_j(\vec{x})$, $\forall j \neq k$.

DBS are then given by $y_k(\vec{x}) = y_j(\vec{x})$, s.t.

$$(\vec{\omega}_k - \vec{\omega}_j)^T \cdot \vec{x} + (\omega_{k,0} - \omega_{j,0}) = 0$$

Now, consider $\vec{x}_A, \vec{x}_B \in C_k$

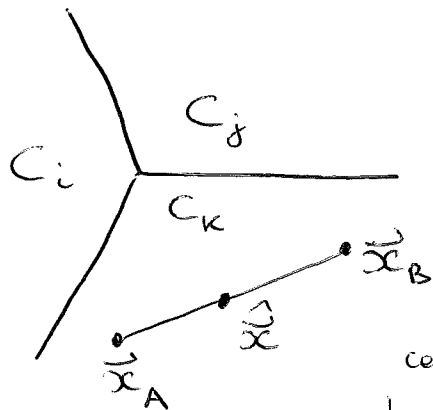line connecting $\vec{x}_A \& \vec{x}_B$

$$\hat{\vec{x}} = \lambda \vec{x}_A + (1-\lambda) \vec{x}_B \qquad 0 \le \lambda \le 1$$
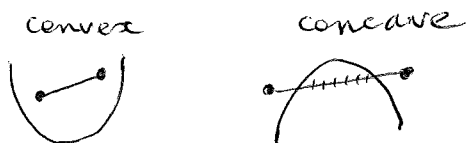
$\Downarrow$ linearity of discriminant f's

$$y_k(\hat{\vec{x}}) = \lambda \underbrace{y_k(\vec{x}_A)}_{\substack{> y_j(\vec{x}_A), \\ \forall j \ne k}} + (1-\lambda) \underbrace{y_k(\vec{x}_B)}_{\substack{> y_j(\vec{x}_B), \\ \forall j \ne k}}$$

$$\underbrace{\phantom{y_k(\hat{\vec{x}}) = \lambda y_k(\vec{x}_A) + (1-\lambda) y_k(\vec{x}_B)}}$$

$$y_k(\hat{\vec{x}}) > y_j(\hat{\vec{x}}) , \quad \forall j \ne k$$

So, $\hat{\vec{x}} \in C_k$ as well.



Since $\vec{x}_A, \vec{x}_B$ are arbitrary, $C_k$ is singly connected & convex, $\forall k$.

# Least squares for classification

Consider a problem with $K$ classes, s.t. $\vec{t}$ are $K$-dim unit vectors.

$$y_k(\tilde{x}) = \vec{w}_k^T \tilde{x} + w_{k,0\#} \implies \tilde{y}(\tilde{x}) = \widetilde{W}^T \tilde{x}$$

$$\widetilde{W} = \overbrace{\left( \begin{array}{c} w_{K,0} \\ w_{k,1} \\ \vdots \\ w_{k,D} \end{array} \right)}^{K \text{ (classes)}} \left. \vphantom{\begin{array}{c} w \\ w \\ \vdots \\ w \end{array}} \right\} D+1 \text{ (\# prms)}$$

or
# entries in $\tilde{x}$

$\underbrace{\phantom{w_{k,D}}}$ $k\underline{th}$ column

$$\widetilde{\tilde{x}} = (1, \vec{x})^T$$
$\phantom{xx}$ "$x_0$

Training set: $\{\vec{x}_n, \vec{t}_n\}$ $\quad n = 1, \ldots, N$

Define $\quad T = \overbrace{\left( t_{n,0} \; t_{n,1} \; \ldots \; t_{n,K-1} \right)}^{K} \left. \vphantom{(} \right\} N$

$\vec{t}_n$ is the $n\underline{th}$ row of $T$

$$\widetilde{X} = \underbrace{\left( \overset{"1}{x}_{n,0} \cdot \ldots \ldots x_{n,D} \right)}_{D+1} \left. \vphantom{(} \right\} N$$

$\widetilde{\tilde{x}}_n$ is the $n\underline{th}$ row of $\widetilde{X}$

$\widetilde{X}\widetilde{W}$ is an $\underset{\substack{\uparrow \\ \#\text{rows}}}{N} \times \underset{\substack{\uparrow \\ \#\text{columns}}}{K}$ matrix like $T$

Then
$$E(\widetilde{W}) = \frac{1}{2} \text{Tr}\left\{ (\widetilde{X}\widetilde{W} - T)^T (\widetilde{X}\widetilde{W} - T) \right\}.$$

Indeed,
$$E(\widetilde{W}) = \frac{1}{2} \sum_{k=1}^{K} \sum_{n=1}^{N} \left( \underbrace{(\widetilde{X}\widetilde{W})_{nk}}_{\displaystyle \sum_{j=0}^{D} \widetilde{x}_{nj} \widetilde{W}_{jk} = \sum_{j=0}^{D} \widetilde{x}_{nj} w_{k,j} =} - t_{nk} \right)^2 \overset{\circ}{=}$$

$$= w_{k,0} + \widetilde{w}_k^T \cdot \widetilde{x}$$

$$\overset{\circ}{=} \frac{1}{2} \sum_{k=1}^{K} \sum_{h=1}^{N} \left( \sum_{j=0}^{D} w_{k,j} \widetilde{x}_{nj} - t_{nk} \right)^2.$$

Then
$$\frac{\partial E}{\partial w_{k,j}} = \frac{1}{2} 2 \sum_{n} \left[ \sum_{j'} w_{k,j'} \widetilde{x}_{nj'} - t_{nk} \right] \widetilde{x}_{nj} = 0,$$
$$\forall k, j$$

$$\underbrace{(\widetilde{X}^T \widetilde{X})}_{(D+1)\times(D+1)} \underbrace{\widetilde{W}}_{(D+1)\times K} = \underbrace{\widetilde{X}^T T}_{(D+1)\times K}$$
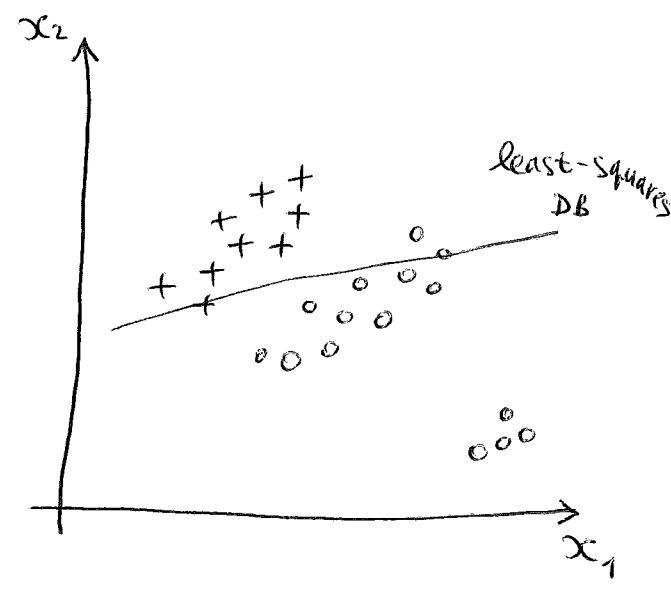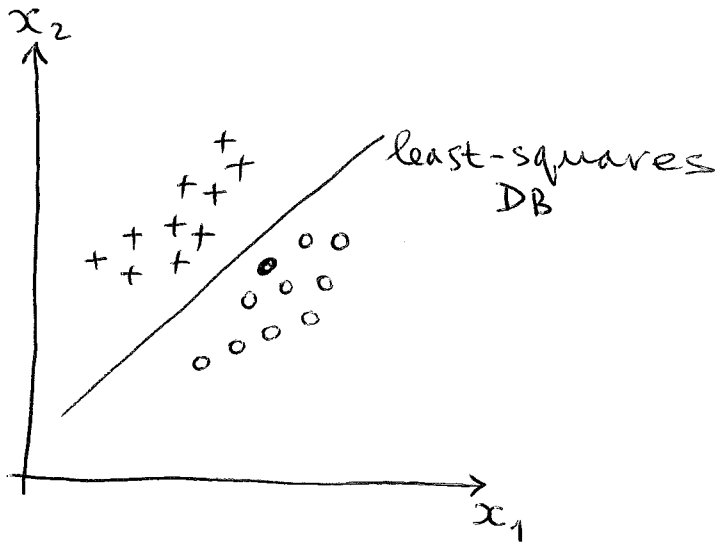$$\underbrace{\phantom{(\widetilde{X}^T \widetilde{X}) \widetilde{W}}}_{(D+1)\times K}$$

Finally, $\quad \widetilde{W} = \underbrace{(\widetilde{X}^T \widetilde{X})^{-1} \widetilde{X}^T}_{\text{"}\widetilde{X}^+ \text{ pseudo-inverse}} T$
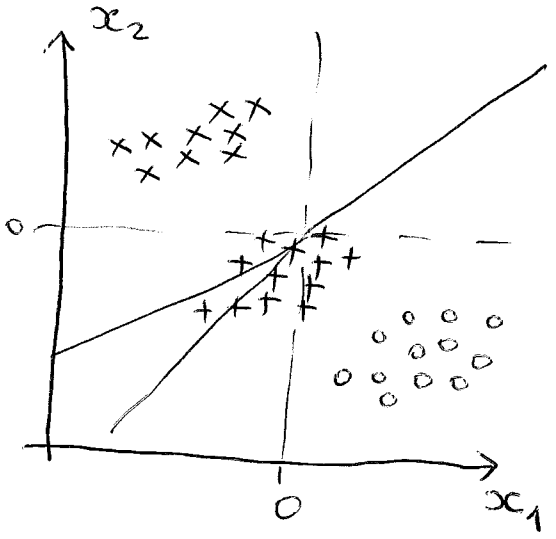$$\text{of } X$$

$$\vec{y}(\widetilde{x}) = \widetilde{W}^T \widetilde{x} = T^T (\widetilde{X}^+)^T \widetilde{x}.$$

This is a closed-form solution which however is sensitive to outliers.

<u>Ex.</u>



And may even fail <u>completely</u>:



poor
prediction

Indeed, least-squares assumes gaussian distr'n of $\vec{t}$'s, and binary target vectors often have non-gaussian distributions.